Table
T =

112

Error tolerance vector

e=[e1, e2, e3, e4, e5, e6, e7]

123

110 Dependency Finder

115 Bayesian Network

120 CaRTSelector

140 CaRTBuilder

130 RowAggregator

Tc
Model-Compressed Table

100

FIGURE 1
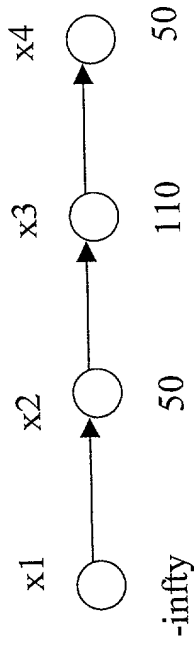
```
procedure Greedy (T(X), ē, G, θ)
Input:     n-attribute table T and n-vector of error tolerances ē;
           Bayesian network G on the set of attributes X and
           threshold  θ on the relative benefit for selecting a
           CaRT predictor.
Output:    A set of materialized (predicted) attributes X_mat (X_pred
           = X - X_mat) and a CaRT predictor for each X_i ∈ X_pred.
begin
1.    X_mat := X_pred := φ
2.    let < X_1, X_2,...,X_n > be the attributes in X sorted in
      topological order of G
3.    for  i :=1,...,n
4.    if π (X_i) = φ then X_mat :=X_mat ∪ {X_i}    /* X_i must be
      materialized if it has no parents in G */
5.    else
6.    M := BuildCaRT (X_mat → X_i, e_i)
7.    if (MaterCost (X_i) / PredCost (X_mat → X_i) >  θ) then X_pred :=
      X_pred ∪ {X_i}
8.    else X_mat := X_mat ∪ {X_i}
9.    end
10.   end
end
```
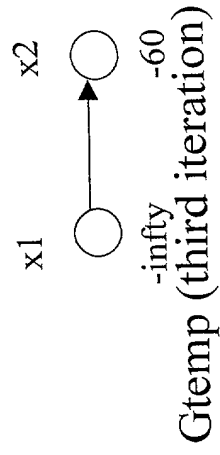
**FIGURE 2:** The Greedy CaRT Selection Algorithm

x1    x2    x3    x4

-infty    50    110    50

(b) Gtemp (first iteration)

## FIGURE 3B

x1    x2    x3    x4

Bayesian Network G

## FIGURE 3A

x1    x2    x4

-infty    -15    45

Gtemp (second iteration)

## FIGURE 3C

x1    x2

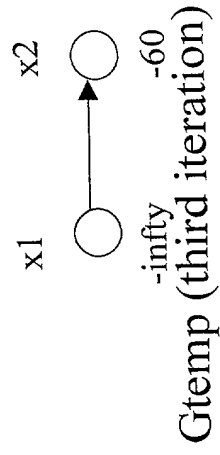-infty    -60

Gtemp (third iteration)

## FIGURE 3D

## FIGUREs 3A-3D

```
procedure MaxIndependentSet (T(X), ē, G, neighborhood() )
Input:      n-attribute table T and n-vector of error tolerances ē;
            Bayesian network G on the set of attributes X and function
            neighborhood()defining the "predictive neighborhood" of a
            node Xᵢ in G (e.g., π (Xᵢ) or β (Xᵢ)).
Output:     A set of materialized (predicted) attributes X_mat (X_pred = X -
            X_mat) and a CaRT predictor PRED (Xᵢ) → Xᵢ for each Xᵢ ∈ X_pred.
begin
1.      X_mat := X, X_pred := ϕ
2.      PRED(Xᵢ) := ϕ for all Xᵢ ∈ X, improve := true
3.      while (improve ≠ false) do
4.              for each Xᵢ ∈ X_mat
5.                      mater_neighbors (Xᵢ) :=
                        (X_mat∩neighborhood(Xᵢ))∪{PRED(X):X ∈ neighborhood
                        (Xᵢ), X ∈ X_pred}-{Xᵢ}
6.                      M := BuildCaRT (Mater_neighbors (Xᵢ)→Xᵢ, eᵢ)
7.                      let PRED (Xᵢ) ⊆ mater_neighbors (Xᵢ) be the set of
                        predictor attributes used in M
8.                      cost_changeᵢ :=0
9.                      for each Xⱼ ∈ X_pred such that Xᵢ ∈ PRED(Xⱼ)
10.                             NEW_PREDᵢ (Xⱼ) := PRED(Xⱼ)-{Xᵢ}∪PRED(Xᵢ)
11.                             M :=BuildCaRT (NEW_PREDᵢ(Xⱼ)→Xⱼ, eⱼ
12.                             set NEW_PREDᵢ (Xⱼ) to the (sub)set of
                                predictor attributes used in M
13.                             cost_changeᵢ := cost_changeᵢ + (PredCost(PRED
                                (Xⱼ)→Xⱼ)-PredCost (NEW_PREDᵢ(Xⱼ)→Xⱼ))
14.                     end
15.             end
16.             build an undirected, node-weighted graph G_temp = (X_mat,
                E_temp) on the current set of materialized
17.             attributes X_mat, where:
18.                     (a) E_temp := { (X,Y) : for all pairs X, Y ∈ X_pred}∪
19.                                           { (Xᵢ,Y): for all Y ∈ X_mat}
20.                     (b) weight (Xᵢ) := MaterCost (Xᵢ) -PredCost (PRED(Xᵢ)
                        → Xᵢ)+cost_changeᵢ for each Xᵢ ∈ X_mat
21.             S := FindWMIS (G_temp)     /* select (approximate) maximum
                weight independent set in G_temp
22.                                             as "maximum-benefit" subset of
                                                predicted attributes  */
23.             if (∑_{X∈S} weight (X) ≤ 0) then improve := false
24.             else/* update X_mat, X_pred, and the chosen CaRT predictors */
25.                     for each Xⱼ ∈ X_pred
26.                             if (PRED(Xⱼ)∩ S = {Xᵢ}) then PRED (Xⱼ) :=
                                NEW_PREDᵢ(Xⱼ)
27.                     end
28.                     X_mat := X_mat - S, X_pred := X_pred ∪ S
29.             end
30.     end /* while */
end
```

**FIGURE 4:**  The MaxIndependentSet CaRT Selection Algorithm

```
procedure LowerBound (N, e_i, b)
Input:     Leaf N for which lower bound on subtree cost is to be
           computed; error tolerance e_i for attribute X_i; bound b
           on the maximum number of internal nodes in subtree
           rooted at N.
Output:    Lower bound L(N) on cost of subtree rooted at N.
begin
1.    for i := to r
2.         minOut [i,0] :=i
3.    for J := 1 to b + 1
4.         minOut[0,j] :=0
5.    l :=0
6.    for i := 1 to r
7.         while x_i - x_{l+1} > 2_{ei}
8.         l :=l = 1
9.     for j := 1 to b + 1
10.        minOut[i,j] := min {minOut[i - 1,j] + 1, minOut [l,j-1]
11.   end
12.   L(N) := ∞
13.   for J := 0 to b
14.        L(N) :=min {L(N), 2j + 1 + j log (|X_i|)+ (j + 1 + minOut
           (r,j+1)) log (|dom(X_i)|)}
15.   L(N) := min {L(N), 2b + 3 + (b + 1) log (|X_i|)+ (b + 2) log
           (|dom(X_i)|)}
16.   return L(N)
end
```

**FIGURE 5:**   Algorithm for Estimating Lower Bound on Subtree Cost
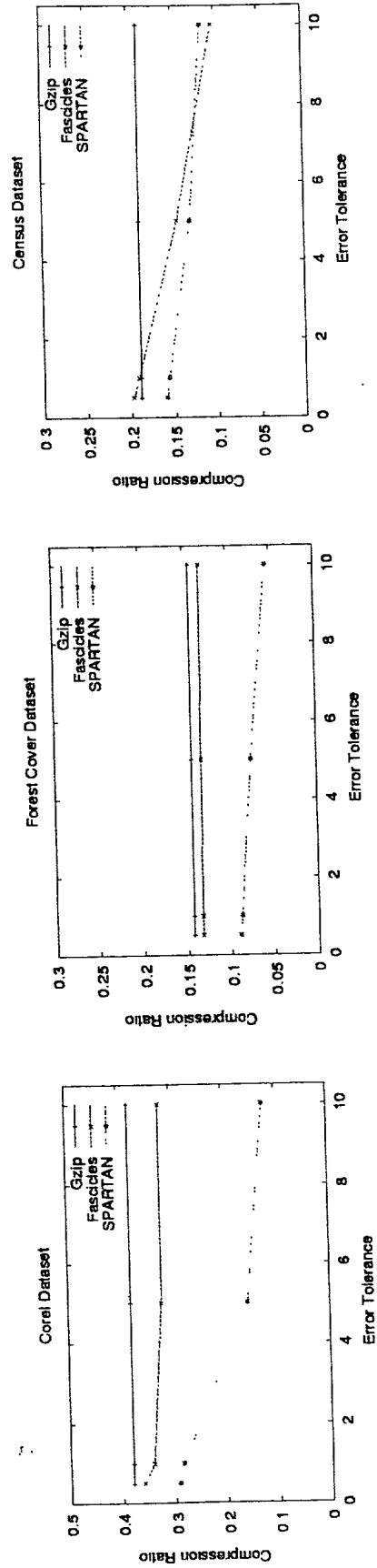
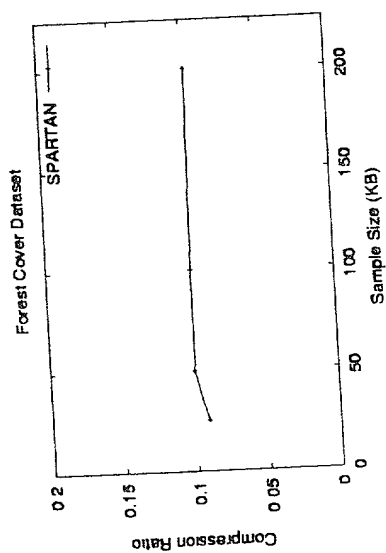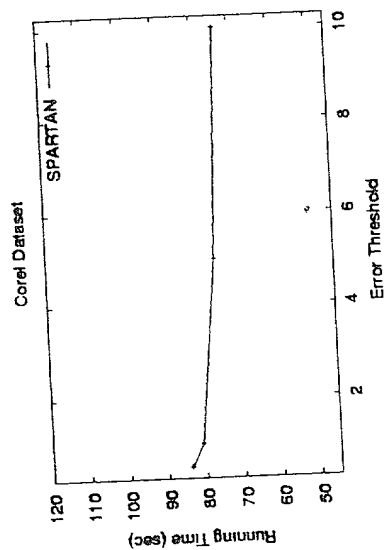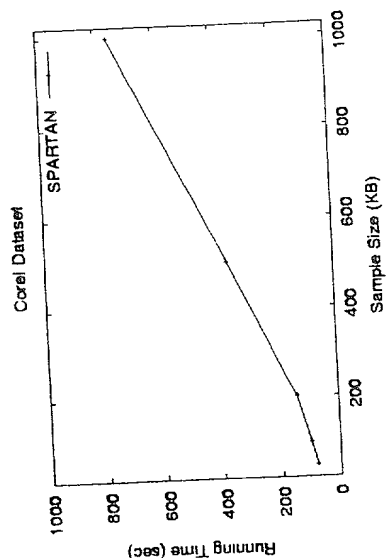Corel Dataset

Forest Cover Dataset

Census Dataset

FIGURE 6

FIGURE 7.